# Genomic Data Access Through BLAST

**https://blast.ncbi.nlm.nih.gov**
Accessing genomic sequence data through BLAST, on the web or using standalone tools
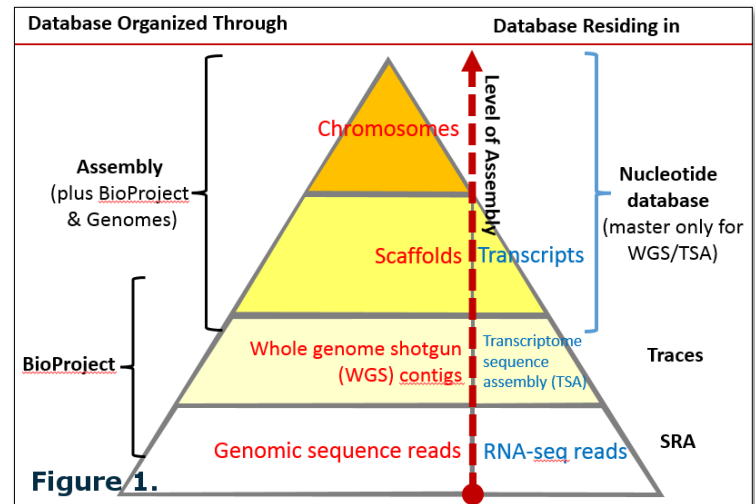National Center for Biotechnology Information • National Library of Medicine • National Institutes of Health • Department of Health and Human Services

## Introduction

Advances in next generation sequencing technology (NGS) have led to the availability of genomic sequence data for an increasing large number of organisms. BLAST searching against these datasets, particularly the annotated assemblies based on raw sequence reads, can provide significant insight into the biology of these biomedically, agriculturally and ecologically important species. However, assembly and annotation from raw sequence reads is a complex process. The availability of best sequence data vary from organism to organism and will require different access strategies. In this document, we will go over the organization of genomic sequence data available from NCBI, and ways to locate the best genomic dataset for the organism of interest for sequence alignment purposes, through both the BLAST homepage at blast.ncbi.nlm.nih.gov/, other web pages/records at NCBI, as well as through standalone tools provided by NCBI.

## Workflow and Organization for Nucleotide Sequence Data

From an NGS-centric point of view, we can organize available nucleotide sequence data based on their volume, degree of assembly, and information density (quality of annotation) as a pyramid structure (Figure 1, right). The figure separates nucleotide sequences into genomic and transcripts entries (left and right), and sorts them by their level of assembly (bottom to top). We should use the records at the top of the pyramid with the highest level of assembly and annotation. It is better to access the data through the organizational databases, e.g., Assembly or BioProject (left half), since they organize and connect related nucleotide records, such as individual chromosomes for a specific organism, into a biologically-relevant collection unit.



**Figure 1.**

## BioProject and Assembly Entries with Genomic Data

A BioProject database record provides a summary of a specific research project and lists all data available from the project. The result below (**A**) is from searching with "Barley[orgn] AND bioproject_assembly[filter]" (**B**). Click the title to open a record (**C**) for more details. The Project Data table lists available nucleotide sequences, with the number linked to actual records (wgs contigs in this case, **D**). The right hand column lists related records in other NCBI databases (insert, **E**). For project with only raw sequence reads, the link will point to Sequence Read Archive (SRA, **F**).

## BioProject and Assembly Entries with Genomic Data (cont.)

An assembly database record provides summary information for a specific genomic assembly. Searching with "barley[orgn]" and filtering for the "Representatives" retrieves a single record (**A**). Its full report (**B**) contains a detailed description of the assembly, with chromosomal level details given in the "Global assembly definition" table (**C**) when available (absent for this specific assembly). The "BLAST the assembly" link in the right hand column (**D**) leads to a Assembly-specific BLAST search form.

## Organism Search Box in the BLAST Homepage

The search box in the BLAST Genome section of the BLAST homepage (blast.ncbi.nlm.nih.gov/) streamlines the process to allow quick access to the best genomic dataset for the input organism. It returns organism-specific nucleotide BLAST pages in the following decreasing level of assembly:

- fully annotated RefSeq chromosomal assembly (go.usa.gov/xptx4)
- Scaffold level of assembly, generally with annotation and some genomic context, but without chromosome-level placement
- WGS level of assembly, with or without annotation, without genomic context or chromosome-level placement
- NT database with the input organism as limit, if all the above fail

The example below locates the genome assembly for the Chinese hamster. It returns the organism-specific BLAST page with annotated RefSeq genome as the target database: type the organism name in the input box to see a suggested list (**E**), select the desired entry from the list (**F**), click "Search" to get to the search page (**G**), and click the "?" icon (**H**) to toggle on a detailed description of the selected database. The example uses the mouse mRNA of the vitamin C synthesis gene (NM_178747.1, **I**) to identify the hamster counterpart.

**Summary** ▾                                                   Send to:

ℹ Filters activated: Latest, Representative, Exclude derived from surveillance project, Exclude anomalous. Clear all to show 19 items.

Assembly for barley cultivar Morex v1.0
Organism: Hordeum vulgare subsp. vulgare (domesticated barley)
Submitter: IPK-Gatersleben                    **A**
Date: 2019/06/02
Assembly level: Scaffold
Genome representation: full
RefSeq category: **representative genome**
GenBank assembly accession: GCA_901482405.1 (**latest**)
RefSeq assembly accession: n/a
Excluded from RefSeq: genome length too large
IDs: 3430831 [UID]   878 [GenBank]     **B**

Full Report ▾                      Send to: ▾      ⬇ **Download Assembly**

**Assembly for barley cultivar Morex v1.0**
Organism name: Hordeum vulgare subsp. vulgare (domesticated barley)
BioSample: SAMEA5598650
BioProject: PRJEB32488
Submitter: IPK-Gatersleben
Date: 2019/06/02
Assembly level: Scaffold
Genome representation: full
RefSeq category: representative genome
Excluded from RefSeq:
- genome length too large
GenBank assembly accession: GCA_901482405.1 (latest)
RefSeq assembly accession: n/a
RefSeq assembly and GenBank assembly identical: n/a
WGS Project: CABEFD01
Genome coverage: 822x
IDs: 3430831 [UID] 11181878 [GenBank]
History (Show revision history)

See Genome Information for Hordeum vulgare

There are 18 assemblies for this organism
See more

go.usa.gov/xptcB

**Access the data**
BLAST the assembly      **D**
Full sequence report
Statistics report
FTP directory for GenBank assembly

**Assembly Information**
Assembly Help
Assembly Basics
NCBI Assembly Data Model

**Related Information**
BioProject
BioSample
Nucleotide INSDC
Sra
Taxonomy
WGS Master

**Global statistics**

| | |
|---|---|
| Total sequence length | 4,833,791,107 |
| Total ungapped length | 4,446,895,020 |
| Gaps between scaffolds | 0 |
| Number of scaffolds | 8 |
| Scaffold N50 | 657,224,000 |
| Scaffold L50 | 4 |
| Number of contigs | 1,030,204 |
| Contig N50 | 19,388 |
| Contig L50 | 55,370 |
| Total number of chromosomes and plasmids | 0 |
| Number of component sequences (WGS or clone) | 8 |

**C**

Assembly Definition | Assembly Statistics

**Global assembly definition**                    Download the full sequence report
The primary assembly unit does not have any assembled chromosomes or linkage groups. Please download the full sequence report for information on the scaffolds.

**BLAST Genomes**

chinese ha     **E**          **Search**

Chinese hamster (taxid:10029)
Chinese hamsters (taxid:10029)
Chinese hairy crab (taxid:95602     **F**
Chinese habu (taxid:103944)
Chinese hawthorn (taxid:510735)
Chinese hare (taxid:112022)

**BLAST Genomes**

Chinese hamster (taxid:10029)          **Search**     **G**
Enter organism common name, scientific name, or tax id.
Human          Mouse          Rat          Microbes

*Cricetulus griseus* (Chinese hamster) Nucleotide BLAST

blastn | blastp | blastx | tblastn | tblastx

BLASTN programs search nucleotide databases using a nucleotide query. more...          Reset page / Bookmark

**Enter Query Sequence**

Enter accession number(s), gi(s), or FASTA sequence(s) ⓘ          Clear     Query subrange ⓘ

NM_178747.1          **I**                                        From
                                                                To

Or, upload file          Choose File   No file chosen          ⓘ
Job Title          Enter a descriptive title for your BLAST search ⓘ

**Choose Search Set**                                   **H**

Database     Genome (CriGri_1.0 reference, Annotation Release 103)  ▾   ⓘ

    Title: Cricetulus griseus CriGri_1.0 [GCF_000223135.1] chromosomes plus unplaced and unlocalized scaffolds (reference assembly in Annotation Release 103)
    Description: The reference assembly set of RefSeq genomic top-level sequences (chromosomes, unplaced and unlocalized scaffolds) in a specific annotation run
    Molecule Type: Genomic
    Update date: 2018/12/18
    Number of sequences: 109152

Exclude
Optional          ☐ Models (XM/XP)
Entrez Query
Optional          Enter an Entrez query to limit search ⓘ

# An Example Web BLAST Search Result



Searching for the vitamin C synthesis gene, using the mouse mRNA as a query and discontiguous megablast to optimize cross-species comparison, finds good matches in the hamster genome.

The graphical overview indicates that all the query is covered by high score matches. The vertical lines indicate exon boundaries (**A**). Detailed alignment are under the "Alignments" tab (**B**). Change the display to "Pairwise with dots for identities," adding translation using "CDS features" checkbox, and sort the hits by "Query start position" (**C**), we will see mismatches highlighted in pink, translation of annotated protein products projected onto the alignment, with the alignments sorted by the natural query position.

## Table 1. Characteristics of Genomic Data Access Through Web BLAST Interface

| Pros | Cons |
|---|---|
| • FAST: distributed computation (splitd)<br>• GUI: Interactive graphical user interface<br>• Extensive links to:<br>   * NCBI records<br>   * Tools (TreeView, Taxonomy Report)<br>• Visually informative format<br>• Coding Sequence Translation<br>• Graphical rendering through SV | • Browser bottleneck<br>• Limited capacity from CPU time restriction<br>• Difficult to<br>   * automate<br>   * incorporate into other workflow<br>• Limited search customization<br>• Limited custom data access |

## Table 2. Alternative Approaches for Genomic Data Access

| Standalone BLAST+ | Vdb-based BLAST |
|---|---|
| • Client-server access through "-remote" using database and computation power at NCBI<br>• Local database access (downloaded from NCBI or formatted locally)<br>• Multi-threaded<br>• Search once and format multiple times through "-outfmt 11" and blast_formatter | • For sequence data stored in vdb format (WGS, TSA & SRA)<br>• Locally stored or on-demand download from prefetch<br>• Multi-threaded<br><br>[both alternatives are command line only with no graphical user interface] |

# Web Access: Some Pros & Cons

BLAST access of genomic sequence data through the web interface has advantages and limitations (Table 1). For genomic BLAST searches requiring high throughput, customization, or workflow integration, alternative approaches may work better. These alternative approaches include the standalone BLAST+ package, the vdb-based BLAST programs from the NCBI sratoolkit, as well as the cloud implementation. Table 2 summarizes some of the characteristics for BLAST+ and vdb-based BLAST tools.

## Locate the Genomic Dataset Using blastdbinfo

The blastdbinfo database provides information about BLAST databases available from NCBI. We can use the einfo tool from the Entrez Program Utilities API to get the basic information about this database (go.usa.gov/xptgr). To locate a database, we can search with text terms through esearch, then refer to the uids (or the **WebEnv** and **QueryKey** value when using **&usehistory=y**) to get details through esummary. For BLAST+ programs, we can use the value from the **<Path>** field as an argument for the **-db** switch in combination with **-remote** to search against that dataset. The example

(right, **A**) searches for the wine grape genomic assembly and parses the <Path> values using the XML parser xtract from EDirect (go.usa.gov/xptgb). The value serves as an argument in a blastn search (**B**).

```
eDirect Commands
esearch -db blastdbinfo -query "wine grape[database organism taxid] AND genomic[blast database type] AND
refseq[blast database source]" | esummary | xtract -pattern DocumentSummary -element Name,Path

OUTPUT
RefSeq Genomic  GP/29760.12992/RefSeq_Genomic
... ...
GCF_000003745.3 genomic/29760/GCF_000003745.3 <<<<<<<<<<Path value we need
```
**A**

**B**
```
blast+ commands
blastn -db genomic/29760/GCF_000003745.3 -remote -query my_query.fa <other switches>
```

As with other Entrez databases, we can use field-limited terms to do more precise searches as shown below.

- genbank[Blast Database Source], gnomon, pdb, refseq, sra, swissprot, trace
- wgs[Blast Sequence Strategy]
- cdna[Blast Sequence Type], genomic, otherdna, protein
- gca_000003225_1[Genome Collection Assembly Name], gcf_000181295_1
- alistipes_inops[Database Organism Taxid], similar to [organism] in other Entrez databases
- 104937[NCBI Genome Project ID]
- aaaa[NCBI WGS Project ID]
- nucleotide[database title]

For sequence data stored in vdb format (WGS, TSA, and SRA), we need to use **blastn_vdb** and **tblastn_vdb** programS from the SRA toolkit to perform searches. For historical rea-

**C**
```
vdb blastn command
blastn_vdb -db "SRR011188 GACC01 AJKK01" -query q -outfmt '6 sseqid length evalue bitscore'
Tabular output
gi|425936258|gb|AJKK01254946.1| 180    9.56e-89       333
gi|425936258|gb|AJKK01254946.1| 125    3.59e-58       231
... ...
gnl|SRA|SRR011188.63012.2      109    2.81e-49       202
gnl|SRA|SRR011188.21834.2      105    4.71e-47       195
gnl|SRA|SRR011188.27702.2      104    1.69e-46       193
... ...
```
**D**
```
taxid2wgs.pl example:
perl taxid2wgs.pl -alias_file "Streptomyces_scabiei_WGS" -title "Streptomyces scabiei WGS"
1930

Content of the vdb database alias file
#
# Alias file created by taxid2wgs.pl/1.0 on Sat Dec 26 11:12:03 2015
TITLE Streptomyces scabiei WGS
VDBLIST JPPX01 JPPW01 LBNJ01
```

sons, blastdbinfo's <Path> value for WGS and TSA entries contains a "*WGS_VDB://*" prefix that must be removed before using the value as argument for the **-db** switch. The example above (**C**) searches three databases (space-separated within quotes) in vdb formats. For WGS datasets with many volumes, NCBI provides a tool, **taxid2wgs.pl**, to facilitate the collection of WGS datasets available for a specific organism using its taxonomy id (taxid). The example collects available WGS datasets for *Streptomyces scabiei* using its taxid, 1930 (**D**). More information on this tool is available at ftp.ncbi.nlm.nih.gov/blast/WGS_TOOLS/

## Download Databases

We can install BLAST databases locally by downloading them from the NCBI ftp site. Table 3 summarizes the available databases, their formats, and FTP subdirectories.

## Other Help Documents

Standalone BLAST+
go.usa.gov/xpt4a

Local vdb_blast:
go.usa.gov/xpt4c

BLAST in the Cloud
github.com/ncbi/blast_plus_docs

**Table 3. Alternative Approaches for Genomic Data Access**

| Databases | Details |
|---|---|
| Standard set (preformatted) | ftp.ncbi.nlm.nih.gov/blast/db/ <br> Download all volumes, e.g., nt.##.tar.gz, Extract and go |
| RefSeq genome assemblies (FASTA) | ftp.ncbi.nlm.nih.gov/genomes/refseq/ <br> Refer to the readme |
| WGS+TSA * (vdb) | ftp.ncbi.nlm.nih.gov/sra/wgs/ <br> Locate project initials using Entrez Nucleotide query: wheat [orgn] AND (wgs_master[prop] OR tsa_master[prop]) |
| WGS + TSA (FASTA) | www.ncbi.nlm.nih.gov/Traces/wgs/ <br> Browse and navigate to individual project to download |
| SRA (vdb) * | ftp.ncbi.nlm.nih.gov/sra/sra-instant/reads/ByRun/sra/ <br> Search in BioProject & follow links to SRA to locate the SRR# |
| *  Using prefetch from the sratoolkit to download datasets in vdb format is strongly recommended!* | |